

COMPUTING THE MINIMUM NUMBER OF HYBRIDISATION
EVENTS FOR A CONSISTENT EVOLUTIONARY HISTORY

M. Bordewich & C. Semple

*Department of Mathematics and Statistics
University of Canterbury
Private Bag 4800
Christchurch, New Zealand*

Report Number: UCDMS2004/21

NOV 2004

COMPUTING THE MINIMUM NUMBER OF HYBRIDISATION EVENTS FOR A CONSISTENT EVOLUTIONARY HISTORY

MAGNUS BORDEWICH AND CHARLES SEMPLE

ABSTRACT. It is now well-documented that the structure of evolutionary relationships between a set of present-day species is not necessarily tree-like. The reason for this is that reticulation events such as hybridisations mean that species are a mixture of genes from different ancestors. Since such events are relatively rare, a fundamental problem for biologists is to determine the smallest number of hybridisation events required to explain a given (input) set of data in a single (hybrid) phylogeny. The main results of this paper show that computing this smallest number is both NP-hard and APX-hard in the case the input is a collection of phylogenetic trees on sets of present-day species. This answers a problem which was raised at a recent conference. As a consequence of these results, we also correct a previously published NP-hardness proof in the case the input is a collection of binary sequences, where each sequence represents the attributes of a particular present-day species. The NP- and APX-hardness of these problems mean that it is unlikely that there is an efficient algorithm for either computing the result exactly, or approximating it to any arbitrary degree of accuracy.

1. INTRODUCTION

Evolutionary trees, also called (rooted) phylogenetic trees, are used in evolutionary biology to represent the ancestral history of a collection of present-day species. However, evolution is not always tree-like because of reticulation events such as hybridisations and lateral gene transfers. Consequently, rooted acyclic digraphs are being used to model reticulate evolution in which there is exactly one vertex that has in-degree zero and where the vertices of out-degree zero represent the present-day species (see, for example, [2, 7, 12, 16]). In such digraphs, vertices with in-degree at least two represent reticulation events. In this paper, we generically call these vertices ‘hybridisation vertices’ and these digraphs ‘hybrid phylogenies’.

Hybridisation events are relatively rare and so a fundamental problem for biologists studying the evolution of species whose past has included hybridisation is the following: given a collection of phylogenetic trees on sets of species that correctly

Date: 19 November 2004.

1991 Mathematics Subject Classification. 05C05; 92D15.

Key words and phrases. Rooted phylogenetic tree, reticulate evolution, hybrid phylogeny, phylogenetic network, agreement forest, rooted subtree prune and regraft.

The first author was supported by the New Zealand Institute of Mathematics and its Applications funded programme *Phylogenetic Genomics* and the second author was supported by the New Zealand Marsden Fund (UOC310).

represent the tree-like evolution of different parts of various species genomes, what is the smallest number of hybridisation events required so that all of the trees in this collection are simultaneously ‘displayed’ by a single hybrid phylogeny. This smallest number sets a lower bound on the degree of hybridisation that has occurred in the evolution of the species under consideration. Posed in this way in [2], this and similar problems have attracted recent interest (see, for example, [6, 7, 18]). The main results of this paper show that not only is computing this smallest number NP-hard, but that computing it is also APX-hard. The latter means that, unless $P=NP$, there is some fixed positive constant c strictly bigger than 1 for which there is no polynomial-time algorithm such that, for all instances, the ratio between the size of the feasible solution outputted by the algorithm and the size of the optimal solution is always smaller than c . In fact, we show that the APX-hardness (and hence the NP-hardness) of computing this smallest number holds even for the simplest case in which the input collection consists of just two phylogenetic trees on the same set of species.

The paper is organised as follows. The next section contains some necessary preliminaries and a mathematical formalisation of the above optimisation problem (which we call MINIMUM HYBRIDISATION). Formal statements of the main results of this paper, as well as a short summary of the complexity classes and concepts used in these results is also included in this section. The proofs of the main results are given in Section 3. Section 4 contains a discussion of the problem *perfect phylogeny with recombination*, previously examined in [7] and [18]. We point out an error in the proof given in [18] that this problem is NP- and APX-hard, and use our earlier results to provide a correct proof. The last section, Section 5, contains some final remarks including some consequences of the work in Section 3 for the computational complexity of computing the ‘rooted subtree prune and regraft distance’ between a pair of phylogenetic trees. In general, the notation and terminology throughout this paper follows [15].

2. PRELIMINARIES AND MAIN RESULTS

For a digraph D and a vertex v of D , we denote the in-degree and out-degree of v by $d^-(v)$ and $d^+(v)$, respectively. A *hybrid phylogeny* or *hybrid* (on X) is an ordered pair $\mathcal{H} = (D; \phi)$ consisting of

- (i) a rooted acyclic digraph D in which the root has out-degree at least two and, for all vertices v with $d^+(v) = 1$, we have $d^-(v) \geq 2$, and
- (ii) a bijective map ϕ from X into the set of vertices of D with out-degree zero.

For completeness, if $|X| = 1$, then the digraph consisting of an isolated vertex v and a map from X into $\{v\}$ is also defined to be a hybrid on X . The set X corresponds to the set of present-day species and is called the *label set* of \mathcal{H} which is denoted by $\mathcal{L}(\mathcal{H})$. Vertices of in-degree at least two (called *hybridisation vertices*) represent hybridisation events and correspond to an exchange of genetic information between

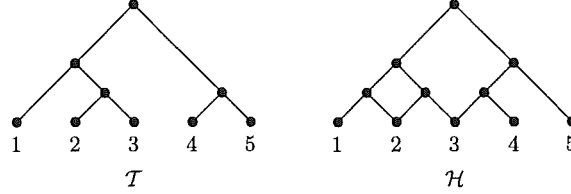


FIGURE 1. A rooted binary phylogenetic tree \mathcal{T} and a hybrid \mathcal{H} displaying \mathcal{T} .

hypothetical ancestors. The *hybridisation number* of \mathcal{H} , denoted $h(\mathcal{H})$, is

$$h(\mathcal{H}) = \sum_{v \neq \rho} (d^-(v) - 1),$$

where ρ denotes the root of \mathcal{H} . Observe that $h(\mathcal{H}) \geq 0$, and $h(\mathcal{H}) = 0$ precisely if D is a rooted tree. Throughout this paper, we adopt the convention that hybrid phylogenies are always drawn with their arcs directed downwards and so omit the arrowheads. A hybrid phylogeny \mathcal{H} with $h(\mathcal{H}) = 2$ is shown in Fig. 1.

A *rooted binary phylogenetic tree* is a special type of hybrid phylogeny in which the root has degree two and all other interior vertices have degree three, and (apart from the root) all vertices have in-degree one.

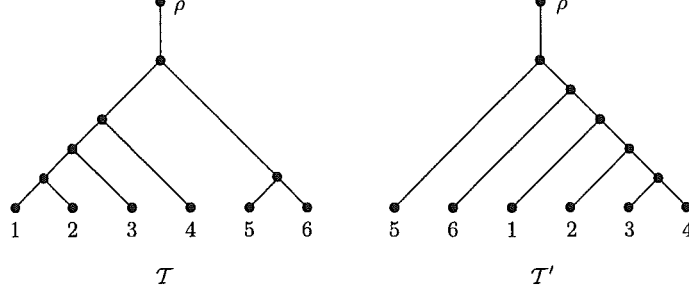
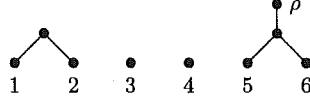
Let \mathcal{T} be a rooted binary phylogenetic X -tree and let \mathcal{H} be a hybrid phylogeny on X . We say that \mathcal{H} *displays* \mathcal{T} if \mathcal{T} can be obtained from a rooted subtree of \mathcal{H} by suppressing degree-two vertices. For example, in Fig. 1, the hybrid \mathcal{H} displays the rooted binary phylogenetic tree \mathcal{T} . For two rooted binary phylogenetic X -trees \mathcal{T} and \mathcal{T}' , we set

$$h(\mathcal{T}, \mathcal{T}') = \min\{h(\mathcal{H}) : \mathcal{H} \text{ is a hybrid on } X \text{ that displays } \mathcal{T} \text{ and } \mathcal{T}'\}.$$

Extending the work in [4], it is shown in [3] that the value $h(\mathcal{T}, \mathcal{T}')$ can be interpreted in terms of a particular type of ‘agreement forest’ for \mathcal{T} and \mathcal{T}' . To make this precise, we need several new definitions.

Let \mathcal{T} be a rooted binary phylogenetic X -tree and let X' be a subset of X . The minimal rooted subtree of \mathcal{T} that connects the vertices of \mathcal{T} labelled by the elements of X' is denoted by $\mathcal{T}(X')$. Furthermore, the *restriction* of \mathcal{T} to X' , denoted $\mathcal{T}|X'$, is the rooted binary phylogenetic tree that is obtained from $\mathcal{T}(X')$ by suppressing any non-root vertices of degree two.

Let \mathcal{T} and \mathcal{T}' be two rooted binary phylogenetic X -trees. For the purposes of the definition of an agreement forest, we regard the root of both \mathcal{T} and \mathcal{T}' as a vertex ρ at the end of a pendant edge adjoined to the original root. Furthermore, we also regard ρ as part of the label sets of \mathcal{T} and \mathcal{T}' , thus we view both label sets as $X \cup \{\rho\}$. An *agreement forest* for \mathcal{T} and \mathcal{T}' is a collection $\{\mathcal{T}_\rho, \mathcal{T}_1, \mathcal{T}_2, \dots, \mathcal{T}_k\}$, where \mathcal{T}_ρ is a rooted tree whose label set \mathcal{L}_ρ includes ρ and $\mathcal{T}_1, \mathcal{T}_2, \dots, \mathcal{T}_k$ are rooted binary phylogenetic trees with label sets $\mathcal{L}_1, \mathcal{L}_2, \dots, \mathcal{L}_k$ such that the following properties are satisfied:

FIGURE 2. Two rooted binary phylogenetic trees T and T' .FIGURE 3. A maximum-good-agreement forest for T and T' .

- (i) The label sets $\mathcal{L}_\rho, \mathcal{L}_1, \mathcal{L}_2, \dots, \mathcal{L}_k$ partition $X \cup \{\rho\}$.
- (ii) For all $i \in \{\rho, 1, 2, \dots, k\}$, $T_i \cong T|_{\mathcal{L}_i} \cong T'|_{\mathcal{L}_i}$.
- (iii) The trees in $\{T(\mathcal{L}_i) : i \in \{\rho, 1, 2, \dots, k\}\}$ and $\{T'(\mathcal{L}_i) : i \in \{\rho, 1, 2, \dots, k\}\}$ are vertex disjoint rooted subtrees of T and T' , respectively.

It is easily seen that if \mathcal{F} is an agreement forest for T and T' , then, up to suppressing non-root vertices of degree two, \mathcal{F} can be obtained from each of T and T' by deleting $|\mathcal{F}| - 1$ edges. An agreement forest for T and T' is a *maximum-agreement forest* if, amongst all agreement forests for T and T' , it has the smallest number of components.

Let $\mathcal{F} = \{T_\rho, T_1, T_2, \dots, T_k\}$ be an agreement forest for T and T' . Let $G_{\mathcal{F}}$ be the directed graph whose vertex set is \mathcal{F} and for which (T_i, T_j) is an arc precisely if $i \neq j$ and either

- (I) the root of $T(\mathcal{L}_i)$ is an ancestor of the root of $T(\mathcal{L}_j)$, or
- (II) the root of $T'(\mathcal{L}_i)$ is an ancestor of the root of $T'(\mathcal{L}_j)$.

Since \mathcal{F} is an agreement forest, the roots of $T(\mathcal{L}_i)$ and $T(\mathcal{L}_j)$, and the roots of $T'(\mathcal{L}_i)$ and $T'(\mathcal{L}_j)$ are not the same. We say that \mathcal{F} is a *good-agreement forest* if $G_{\mathcal{F}}$ is acyclic. Furthermore, if \mathcal{F} contains the smallest number of components over all good-agreement forests for T and T' , we say that \mathcal{F} is a *maximum-good-agreement forest* for T and T' , in which case we denote this value of k by $m_g(T, T')$. Observe that $m_g(T, T') = 0$ if and only if, up to isomorphism, T and T' are identical. To illustrate these definitions, Fig. 3 shows a maximum-good-agreement forest for the two rooted binary phylogenetic trees shown in Fig. 2, where we have adjoined to the root of each of T and T' a pendant edge as described above. The following theorem is established in [3].

Theorem 2.1. *Let T and T' be two rooted binary phylogenetic X -trees. Then*

$$h(T, T') = m_g(T, T').$$

Theorem 2.1 is crucial to the results in this paper. Moreover, because of this theorem, we formally state the optimisation problem MINIMUM HYBRIDISATION as follows.

MINIMUM HYBRIDISATION

Instance: A finite set X , and two rooted binary phylogenetic X -trees T and T' .

Goal: Find a maximum-good-agreement forest \mathcal{F} for T and T' .

Measure: The number of components in \mathcal{F} minus one.

The main results of this paper are Theorem 2.2 and Corollary 2.3.

Theorem 2.2. *The optimisation problem MINIMUM HYBRIDISATION is APX-hard. In particular, there is no polynomial-time approximation scheme for MINIMUM HYBRIDISATION unless $P=NP$.*

It is an immediate consequence of Theorem 2.2 that MINIMUM HYBRIDISATION is NP-hard. Of course, it also means that the (general) fundamental problem described in the introduction is NP- and APX-hard.

Corollary 2.3. *Unless $P=NP$, there is no polynomial-time approximation algorithm for MINIMUM HYBRIDISATION with an approximation ratio better than $\frac{4561}{4560}$.*

We end this section with a short summary of the complexity classes and concepts described in Theorem 2.2 and Corollary 2.3. For further details, we refer the reader to [13] and [1].

For optimisation problems that are NP-hard, an important consideration is the possibility of polynomial-time approximation algorithms. In such an algorithm, one would like to guarantee for all instances that the ratio between the size of the feasible solution outputted by the algorithm and the size of an optimal solution is always smaller than some fixed constant. To treat minimisation and maximisation problems in the same way, we will assume that this ratio is always at least 1. The existence of polynomial-time approximation algorithms varies greatly amongst NP-hard problems. Indeed, there are some NP-hard problems π for which regardless of the size of this fixed constant, there is always such an algorithm. In this case, π is said to exhibit a *polynomial-time approximation scheme* (PTAS). Such problems include the problem of finding a maximum independent set in a planar graph. But then there are other NP-hard problems, such as the (general) travelling salesman problem, for which there exists no polynomial-time approximation algorithm (no matter how big the fixed constant is) unless $P=NP$.

The class APX (also known as MAX SNP) is the class of optimisation problems for which there exists a polynomial-time approximation algorithm for some constant approximation ratio. Within this class, is the class of APX-complete problems. If

an optimisation problem is APX-complete, then it has no polynomial-time approximation scheme unless $P=NP$. Assuming that $P \neq NP$, this implies that there is some fixed constant strictly bigger than 1 for which there is no polynomial-time approximation algorithm such that, for all instances, the approximation ratio is always smaller. To show that an optimisation problem π_2 is APX-hard, it suffices to find an APX-complete problem π_1 and show that there is an ‘ L -reduction’ from π_1 to π_2 . Introduced by Papadimitriou and Yannakakis [13], the reason that this suffices is that L -reductions preserve approximability.

Let π_1 and π_2 be two optimisation problems. An L -reduction from π_1 to π_2 is a pair of polynomial-time computable functions f and g , and a pair of positive constants α and β that satisfy the following properties:

- (i) If I is an instance of π_1 , then $f(I)$ is an instance of π_2 with

$$\text{opt}(f(I)) \leq \alpha \text{opt}(I),$$

where $\text{opt}(I)$ and $\text{opt}(f(I))$ denotes the ‘cost’ of an optimal solution to I and $f(I)$, respectively.

- (ii) If S is a feasible solution of $f(I)$, then $g(S)$ is a feasible solution of I with

$$|\text{opt}(I) - c(g(S))| \leq \beta |\text{opt}(f(I)) - c(S)|,$$

where $c(g(S))$ and $c(S)$ is the ‘cost’ of $g(S)$ and S , respectively.

We establish Theorem 2.2 by using an L -reduction from the APX-complete problem MAXIMUM BOUNDED 3-DIMENSIONAL MATCHING.

3. PROOFS OF THEOREM 2.2 AND COROLLARY 2.3

In this section, we establish Theorem 2.2 and Corollary 2.3. Consider the following problem

MAXIMUM BOUNDED 3-DIMENSIONAL MATCHING (MAX-3DM-B)

Instance: Three disjoint sets X , Y , and Z . A subset Q of $X \times Y \times Z$ of ordered triples with the property that any element of $X \cup Y \cup Z$ appears in at most B triples of Q .

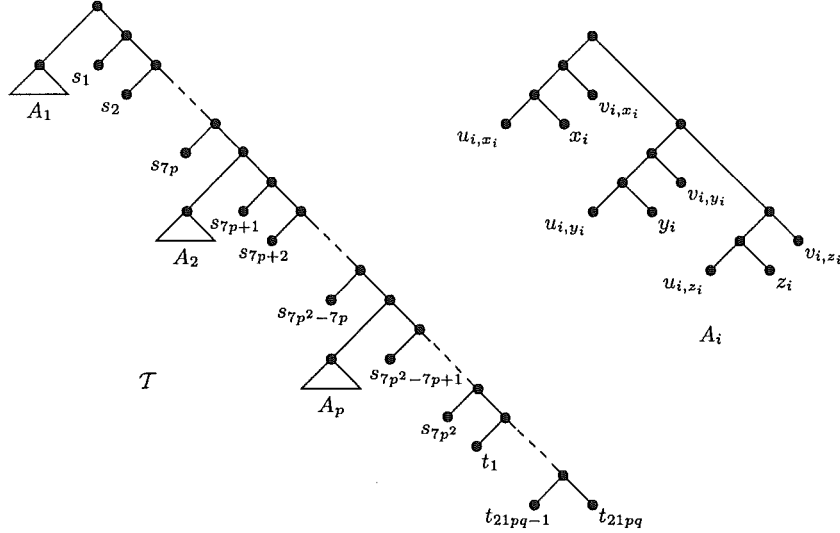
Goal: Find a maximum-sized subset M of Q with the property that no two triples of M agree in any coordinate.

Measure: The cardinality of M .

Kann [10] showed that MAX-3DM-B is APX-complete for all $B \geq 3$. To show that MINIMUM HYBRIDISATION is APX-hard, we show that there is an L -reduction from MAX-3DM-3 to MINIMUM HYBRIDISATION.

Let X , Y , Z , and Q be an instance I of MAX-3DM-3. Let $|Q| = p$. Without loss of generality, we may assume that

$$|X| = |Y| = |Z|$$


 FIGURE 4. The tree \mathcal{T} and its subtrees A_i .

as we can always pad X , Y , or Z with additional elements none of which appear in an ordered triple of Q . Furthermore, we may also assume that

$$q \leq p \leq 3q,$$

where $q = |X| = |Y| = |Z|$. To see this last assumption, we can always remove elements from X , Y , or Z appropriately so that $q \leq p$. Moreover, as each element of $X \cup Y \cup Z$ appears in at most three triples, $p \leq 3q$.

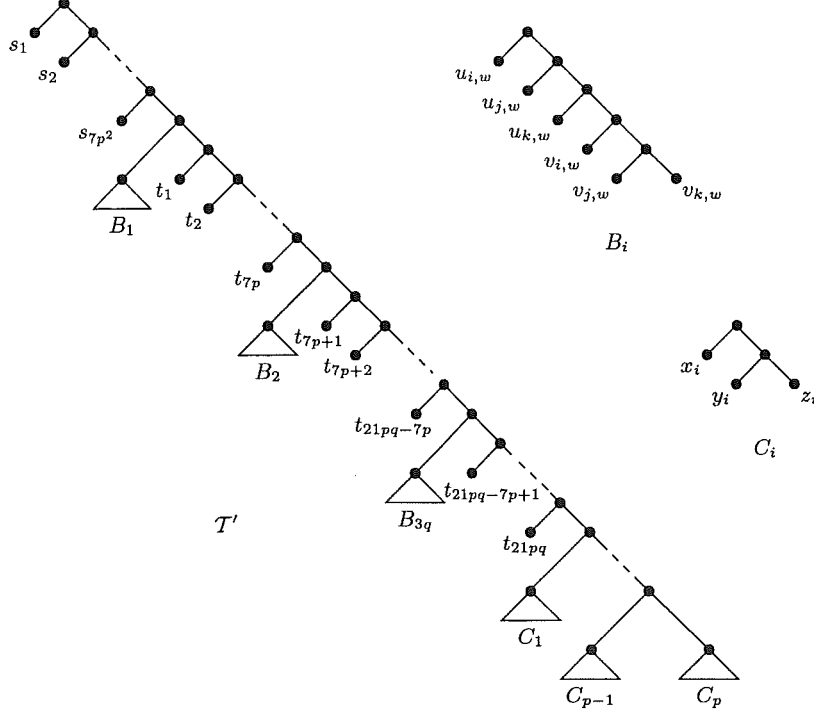
Using the above instance of MAX-3DM-3, we now construct two rooted binary phylogenetic trees \mathcal{T} and \mathcal{T}' with the same label sets. With some modifications, this construction follows the same construction as that used in [9] and [4] to show that a certain related problem is NP-hard but with MAX-3DM-3 replacing EXACT COVER BY 3-SETS (see Section 5 for further details).

Let $Q = \{(x_1, y_1, z_1), (x_2, y_2, z_2), \dots, (x_p, y_p, z_p)\}$. The tree \mathcal{T} is shown in Fig. 4. Each subtree A_i corresponds to exactly one of the triples in Q . The tree \mathcal{T}' is shown in Fig. 5. Each subtree B_i corresponds to an element w of $X \cup Y \cup Z$. The size of the label set of B_i depends on the number of occurrences of w as a coordinate in a triple of Q . Thus it is possible that the label set of B_i is empty which effectively means that there is no such subtree.

The following lemma is central to the proof that MINIMUM HYBRIDISATION is APX-hard.

Lemma 3.1. *Q contains a 3-dimensional matching of size k if and only if there is a good-agreement forest for \mathcal{T} and \mathcal{T}' of size*

$$1 + 6k + 7(p - k) = 7p - k + 1.$$

FIGURE 5. The tree T' , and its subtrees B_i and C_i .

In particular, $h(T, T') = 7p - \text{opt}(Q)$.

Proof. Suppose Q contains a 3-dimensional matching M of size k . We can obtain a good-agreement forest \mathcal{F}_M of size $7p - k + 1$ for T and T' by making the following edge deletions to T and then suppressing any resulting non-root degree-two vertices:

- (i) For each i , delete the edge attaching A_i to the rest of T .
- (ii) For each i , if A_i corresponds to an ordered triple in M , delete each of the pendant edges attaching x_i , y_i , and z_i , and then delete each of the edges attaching the subtrees containing u_{i,x_i} and v_{i,x_i} , u_{i,y_i} and v_{i,y_i} , u_{i,z_i} and v_{i,z_i} . Thus, in this case, each A_i is broken into 6 components.
- (iii) For each i , if A_i does not correspond to an ordered triple in M , then delete each of the pendant edges attaching u_{i,x_i} , v_{i,x_i} , u_{i,y_i} , v_{i,y_i} , u_{i,z_i} , and v_{i,z_i} . In this case, each A_i is broken into 7 components.

Using the fact that M is a 3-dimensional matching and that each B_i corresponds to a particular element of $X \cup Y \cup Z$, it is straightforward to deduce that \mathcal{F}_M can also be obtained from T' by deleting appropriate edges. Hence \mathcal{F}_M is an agreement forest for T and T' . A routine check now shows that \mathcal{F}_M is also a good-agreement forest.

To prove the converse, let $S = \{s_1, s_2, \dots, s_{7p^2}, t_1, t_2, \dots, t_{21pq}\}$. Let \mathcal{F} be an agreement forest for \mathcal{T} and \mathcal{T}' of size at most $7p+1$. We first show that if T_j is a tree in \mathcal{F} with label set $\mathcal{L}(T_j)$, then if $\mathcal{L}(T_j) \cap \mathcal{L}(A_i) \neq \emptyset$ it follows that $\mathcal{L}(T_j) \subseteq \mathcal{L}(A_i)$, and if $\mathcal{L}(T_j) \cap \mathcal{L}(B_i) \neq \emptyset$ it follows that $\mathcal{L}(T_j) \subseteq \mathcal{L}(B_i)$.

Let T_j be a tree in \mathcal{F} , and first assume that for some i the set $\mathcal{L}(T_j) \cap \mathcal{L}(A_i)$ is non-empty and contains at least one element x of $\mathcal{L}(T)$ not in $\mathcal{L}(A_i)$. Suppose that $x \in \mathcal{L}(A_{i'})$ for some $i' \neq i$. Then, since \mathcal{F} is an agreement forest for \mathcal{T} and \mathcal{T}' , there are at least $7p$ members of S that appear as singletons in \mathcal{F} (those in the chain between A_i and $A_{i'}$). Since the label set of no component in \mathcal{F} contains the entire label set of A_i (for any i), this implies that \mathcal{F} contains at least $7p+p$ components; a contradiction. Now suppose that $x \in S$. Then, since \mathcal{F} is an agreement forest for \mathcal{T} and \mathcal{T}' , we have $|\mathcal{L}(T_j) \cap \{s_{7p^2-7p+1}, \dots, s_{7p^2}\}| \leq 1$. Hence at least $7p-1$ of these labels appear as singletons in \mathcal{F} , again leading to a contradiction. Effectively, this means that to obtain \mathcal{F} from \mathcal{T} each edge joining an A_i to the rest of \mathcal{T} is deleted. The result for $\mathcal{L}(T_j) \cap \mathcal{L}(B_i) \neq \emptyset$ follows by similar reasoning.

Now suppose that \mathcal{F} is a (good) agreement forest of size $1 + 6k + 7(p-k) = 7p - k + 1$. Fixing i , consider A_i . By the argument above, there is a subset of the components of \mathcal{F} in which the union of the label sets is the label set of A_i . Since no component can contain labels from more than one $B_{i'}$, a routine check shows that this subset must have at least 6 elements and, moreover, this subset has exactly 6 elements only if the partition of $\mathcal{L}(A_i)$ induced by the label sets is

$$\{\{x_i\}, \{y_i\}, \{z_i\}, \{u_{i,x_i}, v_{i,x_i}\}, \{u_{i,y_i}, v_{i,y_i}\}, \{u_{i,z_i}, v_{i,z_i}\}\}.$$

It now follows that each A_i contributes at least 6 components to \mathcal{F} . An important observation at this point is that regardless of the composition of \mathcal{F} , it is always a good-agreement forest.

Since \mathcal{F} has $1 + 6k + 7(p-k)$ components, it follows from the last paragraph that at least k of the A_i 's are 'partitioned' into 6 parts as described above. Let A_i and A_j be two such subtrees, and consider the associated triples (x_i, y_i, z_i) and (x_j, y_j, z_j) . Suppose that one of the components agree. Without loss of generality, we may assume that $x_i = x_j$. Since A_i and A_j are both partitioned into 6 parts, $\{u_{i,x_i}, v_{i,x_i}\}$ is the label set of one component of \mathcal{F} and $\{u_{j,x_j}, v_{j,x_j}\}$ is the label set of another component of \mathcal{F} . But then, in \mathcal{T}' , the minimal subtree connecting u_{i,x_i} and v_{i,x_i} and the minimal subtree connecting u_{j,x_j} and v_{j,x_j} are not disjoint; a contradiction. Thus (x_i, y_i, z_i) and (x_j, y_j, z_j) have no coordinates in common. We conclude that Q contains a 3-dimensional matching of size k . This completes the proof of the lemma. \square

Theorem 2.2 is restated here for convenience.

Theorem 2.2. *The optimisation problem MINIMUM HYBRIDISATION is APX-hard. In particular, unless $P=NP$, there is no polynomial-time approximation scheme for MINIMUM HYBRIDISATION.*

Proof. To establish the result it suffices to show that there is an L -reduction from MAX-3DM-3 to MINIMUM HYBRIDISATION. First note that by picking any triple m

in Q and removing all other triples which agree with m in at least one coordinate (thus removing at most 7 triples including the one originally picked), and then picking another triple from the resulting set and continuing this process, we observe that $\text{opt}(Q) \geq \frac{p}{7}$; that is

$$(1) \quad p \leq 7 \text{opt}(Q).$$

Let I be an instance of MAX-3DM-3, and let $f(I)$ be the function that maps I to \mathcal{T} and \mathcal{T}' , an instance of MINIMUM HYBRIDISATION as described prior to Lemma 3.1. Clearly, this mapping is computable polynomial time in the size of I . Furthermore, by Lemma 3.1 and (1),

$$\begin{aligned} h(\mathcal{T}, \mathcal{T}') &= 7p - \text{opt}(Q) \\ &\leq 7(7 \text{opt}(Q)) - \text{opt}(Q) \\ &= 48 \text{opt}(Q). \end{aligned}$$

It now follows that (i) in the definition of an L -reduction holds with $\alpha = 48$.

To see that (ii) holds, let \mathcal{F} be an agreement forest for \mathcal{T} and \mathcal{T}' of size $S_2 + 1 = 7p - k + 1$. Let g be the function that maps \mathcal{F} to the feasible solution of I of size $S_1 = k$ as described at the end of the proof of Lemma 3.1. Again, g can be computed in polynomial time. Then $S_2 = 7p - S_1$, and so

$$\begin{aligned} 7p - \text{opt}(Q) &= h(\mathcal{T}, \mathcal{T}') \\ \Leftrightarrow 7p - S_1 - (7p - \text{opt}(Q)) &= S_2 - h(\mathcal{T}, \mathcal{T}') \\ \Leftrightarrow \text{opt}(Q) - S_1 &= S_2 - h(\mathcal{T}, \mathcal{T}'). \end{aligned}$$

It now follows that (ii) in the definition of an L -reduction also holds with $\beta = 1$. This completes the proof of the theorem. \square

Chlebík and Chlebíková [5] recently showed that, unless $P=NP$, there is no polynomial-time approximation algorithm for MAX-3DM-3 with an approximation ratio better than $\frac{95}{94}$. Using the L -reduction in the proof of Theorem 2.2 and, in particular, the values $\alpha = 48$ and $\beta = 1$, we get Corollary 2.3.

Corollary 2.3. *Unless $P=NP$, there is no polynomial-time approximation algorithm for MINIMUM HYBRIDISATION with an approximation ratio better than $\frac{4561}{4560}$.*

Proof. Suppose that there is such an algorithm and suppose that $P \neq NP$. Then using the notation and terminology in the proof of Theorem 2.2, we have

$$\begin{aligned} \frac{S_2}{h(\mathcal{T}, \mathcal{T}')} &< \frac{4561}{4560} \\ \Leftrightarrow \frac{S_2 - h(\mathcal{T}, \mathcal{T}')}{h(\mathcal{T}, \mathcal{T}')} &< \frac{4561}{4560} - 1 = \frac{1}{4560} \end{aligned}$$

But $h(\mathcal{T}, \mathcal{T}') \leq 48 \text{opt}(Q)$, and so

$$\frac{1}{48 \text{opt}(Q)} \leq \frac{1}{h(\mathcal{T}, \mathcal{T}')}.$$

Furthermore, $S_2 - h(\mathcal{T}, \mathcal{T}') = \text{opt}(Q) - S_1$. Therefore

$$\begin{aligned} & \frac{1}{48 \text{opt}(Q)} (\text{opt}(Q) - S_1) < \frac{1}{4560} \\ \Leftrightarrow & 1 - \frac{48}{4560} < \frac{S_1}{\text{opt}(Q)} \\ \Leftrightarrow & \frac{94}{95} < \frac{S_1}{\text{opt}(Q)}. \end{aligned}$$

This last inequality implies that MAX-3DM-3 has a polynomial-time approximation algorithm with an approximation ratio better than $\frac{95}{94}$, contradicting Chlebík and Chlebíkova's result. This completes the proof of the corollary. \square

4. PERFECT PHYLOGENETIC NETWORKS WITH RECOMBINATION

Perfect phylogenetic network with recombination is a problem that has a very similar flavour to that of MINIMUM HYBRIDISATION, and has been studied by Gusfield *et al.* [7] and Wang *et al.* [18]. Like MINIMUM HYBRIDISATION, the goal of this problem is to compute the minimum number of hybridisation events that is required to explain a given input, where in this case the input is a collection of binary sequences. It is shown in [18] that perfect phylogeny with recombination is NP- and APX-hard, however, an assertion in the NP-hardness proof is incorrect. In terms of the language used in this paper, this assertion states that if the 'rooted subtree prune and regraft distance' (see Section 5 for definition) of two rooted binary phylogenetic trees is k , then there is a hybrid phylogeny with k hybridisation vertices each of in-degree two that displays both trees. In [3], explicit examples are given to show that this does not always hold. In this section, we verify the NP- and APX-hardness of the perfect phylogenetic network with recombination problem using the hardness results of MINIMUM HYBRIDISATION.

Although perfect phylogenetic network with recombination could be stated in terms of hybrid phylogenies, we formally state the problem in the language given in [7, 18]. An (n, m) -*phylogenetic network* \mathcal{N} is a rooted acyclic digraph with exactly n vertices of out-degree zero in which each vertex other than the root has either one or two incoming edges, and each vertex of \mathcal{N} is labelled with a binary sequence of length m . A vertex with two incoming edges is called a *recombination vertex*. Each integer in $\{1, 2, \dots, m\}$ is assigned to exactly one edge of \mathcal{N} that is not directed towards a recombination vertex. Beginning with the root which is labelled with the all-0 sequence, each of the binary sequences labelling the other vertices is based on the binary sequence of its parent and the incoming edge (in the case it is a non-recombination vertex) or its parents (in the case it is a recombination vertex). In particular, the sequences satisfy the following properties:

- (I) If v is a non-recombination vertex with incoming edge e , then the sequence labelling v is obtained from the sequence labelling its parent by changing the i -th element from 0 to 1 for each integer i assigned to e . If no integer is assigned to e , then the sequence labelling v is the same as its parent.

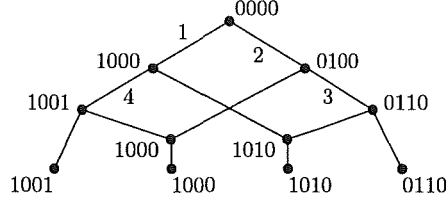


FIGURE 6. A phylogenetic network.

- (II) If v is a recombination vertex, then, for some positive integer p strictly between 1 and m , the sequence labelling v is the concatenation of the first p elements of the sequence labelling one of its parents and the last $m - p$ elements of its other parent.

As an example, a phylogenetic network is shown in Fig. 6. For each recombination vertex in this example, the first two elements in the associated sequence come from its ‘left’ parent and the second two elements come from its ‘right’ parent.

Let B be a collection of n binary sequences of length m . An (n, m) -phylogenetic network \mathcal{N} *explains* B if the n vertices of out-degree zero are bijectively labelled with the elements of B . For example, the phylogenetic network in Fig. 6 explains the collection $\{1001, 1000, 1010, 0110\}$ of binary sequences.

Over all phylogenetic networks that explain B , we are interested in finding one with the smallest number of recombination vertices. We denote this smallest number by $r(B)$. The perfect phylogenetic network with recombination problem is formally stated as follows:

PERFECT PHYLOGENY WITH RECOMBINATION

Instance: A set B of n binary sequences of length m .

Goal: Find a (n, m) -phylogenetic network \mathcal{N} that explains B with the minimum number of recombination vertices.

Measure: The number of recombination vertices in \mathcal{N} .

The motivation for PERFECT PHYLOGENY WITH RECOMBINATION is similar to that for MINIMUM HYBRIDISATION except that, rather than having an input collection consisting of rooted binary phylogenetic trees, we now have an input collection consisting of binary sequences. Each sequence represents a present-day species and, in such a sequence, each coordinate represents some attribute (or character) of the species. A 1 usually indicates that the species under consideration has this particular attribute, while a 0 indicates that the species does not have this attribute. Observe that $0 \rightarrow 1$ is the only allowable transition. The reason for the wording “perfect phylogeny” is that the classical perfect phylogeny problem can be interpreted as the problem of deciding if there is a phylogenetic network with no recombination vertices that explains B .

As mentioned at the beginning of this section, the proof in [18] that establishes the NP-hardness of PERFECT PHYLOGENY WITH RECOMBINATION uses an incorrect assertion. However, the result itself is correct as we next show.

To prove the NP-hardness of PERFECT PHYLOGENY WITH RECOMBINATION, we use a reduction from MINIMUM HYBRIDISATION. We remark here that, even if the NP-hardness proof in [18] was correct, it appears that there is no simple reduction from PERFECT PHYLOGENY WITH RECOMBINATION to MINIMUM HYBRIDISATION. Let T and T' be two rooted binary phylogenetic X -trees, where $|X| = n$. For T and T' , bijectively label the edges with the elements of $\mathcal{C} = \{\chi_1, \chi_2, \dots, \chi_{2(n-1)}\}$ and $\mathcal{C}' = \{\chi'_1, \chi'_2, \dots, \chi'_{2(n-1)}\}$, respectively. Note that both T and T' have $2(n-1)$ edges. Each of the elements in \mathcal{C} and \mathcal{C}' represent a binary character with states 0 and 1. For each vertex v and v' of T and T' , respectively, we associate the binary sequence in which the i -th element is 1 if and only if χ_i (resp. χ'_i) labels an edge on the path from v to the root of T (resp. T'). For each x in X , concatenate the sequence labelling x in T with the sequence labelling x in T' . Let B be the resulting collection of n binary sequences of length $4(n-1)$. This construction is the same as that originally used in [18]. Analogous to Lemma 3.1, the following lemma is central to proving the NP-hardness (and APX-hardness) of PERFECT PHYLOGENY WITH RECOMBINATION.

Lemma 4.1. *Let T and T' be two rooted binary phylogenetic X -trees, and let B be a collection of binary sequences that is constructed from T and T' as above. Then*

$$r(B) = h(T, T').$$

Proof. We first show that $r(B) \leq h(T, T')$. Let \mathcal{H} be a hybrid phylogeny on X that displays T and T' , and has the property that $h(\mathcal{H})$ is minimised. Let ρ denote the root of \mathcal{H} . Because of minimality and the fact that we have only two trees, each hybridisation vertex of \mathcal{H} has in-degree two. By deleting and contracting edges if necessary, we may assume that all the edges of \mathcal{H} are used in some simultaneous displaying of T and T' . Furthermore, by refining vertices if necessary, we may also assume that if a vertex in \mathcal{H} has in-degree two, then it has out-degree one. Now colour each vertex and edge of \mathcal{H} green or red depending upon whether it is used by T or T' , respectively, under the simultaneous displaying of T and T' . Every vertex and edge is coloured with at least one colour. We will call a vertex or edge *monochromatic* if it is only coloured with one colour; otherwise we call it *bichromatic*. We force the root of \mathcal{H} to be bichromatic as follows. In the case that the root of one of the trees, T' say, is identified with a non-root vertex of \mathcal{H} , we will colour ρ and the edges of a directed path from ρ to this non-root vertex of \mathcal{H} red, and view this path as part of T' . The reason for this will be made clear soon. We next assign a binary sequence to each vertex of \mathcal{H} based on this colouring.

As in the case of the sequences in B , the labelling comes in two parts. The root ρ is given the all-0 sequence. Consider the restriction of \mathcal{H} to the green vertices and edges. For each green vertex $v \neq \rho$, assign it the first part of the sequence labelling the vertex of T corresponding to v . If v has degree two in this restriction, assign it the labelling of the first vertex ‘above’ it that has degree three or, in the

case this vertex is the root, degree two. Now consider the restriction of \mathcal{H} to the red vertices and edges. For each red vertex $v \neq \rho$, assign it the second part of the sequence labelling the vertex of T' corresponding to v . If v has degree two in this restriction, assign it the labelling of the first vertex 'above' it that has degree three or, in the case this vertex is the root, degree two. After this labelling, all of the bichromatic vertices of \mathcal{H} have been assigned a sequence with both parts. If v is a monochromatic vertex of \mathcal{H} coloured green, then the second part of its sequence label is the same as the second part of the sequence labelling the first bichromatic vertex that is met on the unique green path from v to ρ . Furthermore, if v is a monochromatic vertex of \mathcal{H} coloured red, then the first part of its sequence label is the same as the first part of the sequence labelling the first bichromatic vertex that is met on the unique red path from v to ρ . Since ρ is bichromatic, this is well-defined.

This direction of the proof is completed by showing that \mathcal{H} with this sequence labelling of the vertices is a phylogenetic network \mathcal{N} that explains B . Clearly, there is a one-to-one correspondence between the elements of B and the vertices of \mathcal{N} of out-degree zero. Furthermore, as \mathcal{H} has the property that the out-degree of each hybridisation vertex v is one, and the edges directed into v are different colours and monochromatic, the sequence assigned to v is of the type described in (II) of the definition of a phylogenetic network. Because of the way in which the elements in B are constructed and the way in which the sequences are assigned to the vertices of \mathcal{H} from the sequences labelling the vertices of T and T' , it is now easily seen that \mathcal{N} is a phylogenetic network that explains B . Hence $r(B) \leq h(T, T')$.

To show that $r(B) \geq h(T, T')$, we can use Claim 2 in the second part of the proof of Theorem 1 in [18] which implies that if there is a phylogenetic network \mathcal{N} that explains B and has k recombination vertices, then the underlying rooted acyclic digraph can be modified to give a rooted acyclic digraph that displays T and T' , and has k recombination vertices, where each recombination vertex has in-degree two. In particular, there is a hybrid phylogeny \mathcal{H} on X that displays T and T' with $h(\mathcal{H}) = k$. Thus $r(B) \geq h(T, T')$. \square

The NP-hardness of PERFECT PHYLOGENY WITH RECOMBINATION follows immediately from the next theorem.

Theorem 4.2. *The optimisation problem PERFECT PHYLOGENY WITH RECOMBINATION is APX-hard.*

Proof. Because of the strength of Lemma 4.1, the proof is straightforward. Let T and T' be an instance I of MINIMUM HYBRIDISATION, and let $f(I)$ be the function that maps T and T' to B , an instance of PERFECT PHYLOGENY WITH RECOMBINATION as described prior to Lemma 4.1. Evidently, this mapping takes polynomial time in the size of T and T' . Furthermore, by Lemma 4.1, $r(B) = h(T, T')$ and so (i) in the definition of an L -reduction holds with $\alpha = 1$.

Now let \mathcal{N} be a phylogenetic network that explains B with S_2 recombination vertices. Let g be the function that maps \mathcal{N} to the feasible solution of T and T' of size $S_1 = S_2$ as described in the last paragraph of the proof of Lemma 4.1. Note

that, as detailed in [18], this mapping can be computed in polynomial time. As $r(B) = h(T, T')$, it follows that

$$S_1 - h(T, T') = S_2 - r(B).$$

Thus (ii) in the definition of an L -reduction holds with $\beta = 1$. \square

The proof of Corollary 4.3 is analogous to that used to prove Corollary 2.3. We omit the details.

Corollary 4.3. *Unless $P=NP$, there is no polynomial-time approximation algorithm for PERFECT PHYLOGENY WITH RECOMBINATION with an approximation ratio better than $\frac{4561}{4560}$.*

5. SOME FINAL REMARKS

Historically, one of the main tools for understanding and modelling reticulate evolution is a graph-theoretic operation called ‘rooted subtree prune and regraft’. The reason for this is that a single rooted subtree prune and regraft operation can be used to model a single reticulation event (see [2, 8, 11, 12, 16]). Moreover, for a pair of rooted binary phylogenetic X -trees, the ‘rooted subtree prune and regraft distance’ between the two trees provides a lower bound to $h(T, T')$ (see [3, 17]). It is stated, but not verified, in [9] that computing this distance is APX-hard. In this section, we verify this result as well as show that, unless $P=NP$, there is no polynomial-time approximation algorithm for computing this distance with an approximation ratio better than $\frac{4561}{4560}$. As we will soon see, it is no coincidence that this ratio is the same as that in Corollary 2.3.

Let T be a rooted binary phylogenetic X -tree. As in the definition of an agreement forest, for the purposes of the upcoming definition, we regard the root of T as a vertex ρ at the end of a pendant edge (called the *root edge*) adjoined to the original root. Let $e = \{u, v\}$ be an edge of T that is not the root edge, where u is the vertex that is in the path from the root of T to v . Let T' be the rooted binary phylogenetic tree obtained from T by deleting e and then adjoining a new edge f between v and the component C_u that contains u as follows. Create a new vertex u' which subdivides an edge in C_u , and adjoin f between u' and v , and then suppress the degree-two vertex u . We say that T' has been obtained from T by a *rooted subtree prune and regraft* (rSPR) operation. We define the *rSPR distance* between two arbitrary rooted binary phylogenetic X -trees T and T' to be the minimum number of rooted subtree prune and regraft operations that is required to transform T into T' . This distance is denoted by $d_{\text{rSPR}}(T, T')$. It is well-known that, for any such pair of trees, one can always obtain one from the other by a sequence of single rSPR operations. Thus this distance is well-defined.

Like the value $h(T, T')$, the value $d_{\text{rSPR}}(T, T')$ can be written in terms of agreement forests. Recall that a maximum-agreement forest for T and T' is an agreement forest with the smallest number of components over all agreement forests for T and T' . Let $m(T, T')$ denote the size of such a forest minus one. Note that there is no reference to ‘good’ in this definition. The following theorem is established in [4].

Theorem 5.1. *Let T and T' be two rooted binary phylogenetic X -trees. Then*

$$d_{\text{rSPR}}(T, T') = m(T, T').$$

Analogous to MINIMUM HYBRIDISATION, we formally state the problem of computing the rSPR distance between two arbitrary rooted binary phylogenetic X -trees in terms of agreement forests.

MINIMUM RSPR

Instance: A finite set X , and two rooted binary phylogenetic X -trees T and T' .

Goal: Find a maximum-agreement forest \mathcal{F} for T and T' .

Measure: The number of components in \mathcal{F} minus one.

Originally thought to be proved in [9], the NP-hardness of MINIMUM RSPR is established in [4] using the original reduction from “Exact Cover by 3-Sets (X3C)” and revising the definition of maximum-agreement forest given in [9] to that described in this paper. This reduction takes an instance of X3C and converts it into a pair of rooted binary phylogenetic trees with the same label sets for which the instance has an exact cover if and only if the two trees has an agreement forest of a certain size. The reduction used in the proof of Theorem 2.2 closely follows this original reduction with MAX-3DM-3 replacing the closely related problem X3C. Indeed, because any maximum-agreement forest for the two rooted binary phylogenetic trees T and T' shown in Fig. 4 and 5, respectively, is also a good agreement forest for T and T' , the pair of trees generated by this reduction always has the property that $d_{\text{rSPR}}(T, T') = h(T, T')$. Consequently, the proofs of Theorem 2.2 and Corollary 2.3 also establish Theorem 5.2. This first part of this theorem verifies a result that is stated without proof in [9].

Theorem 5.2. *The optimisation problem MINIMUM RSPR is APX-hard. Furthermore, unless $P=NP$, there is no polynomial-time approximation algorithm for MINIMUM RSPR with an approximation ratio better than $\frac{4561}{4560}$.*

We end this section by considering what approximation ratios can be achieved in polynomial time for MINIMUM HYBRIDISATION and MINIMUM RSPR. Currently, we do not know of any polynomial-time approximation algorithm for MINIMUM HYBRIDISATION. However, Rodrigues *et al.* [14] describe a polynomial-time 3-approximation algorithm for MINIMUM RSPR. Intuitively, this algorithm builds an agreement forest locally. One might hope that this algorithm extends to MINIMUM HYBRIDISATION, but, due to the additional global condition on a good-agreement forest, it seems unlikely that such an approach will work.

REFERENCES

- [1] G. Ausiello, P. Crescenzi, G. Gambosi, V. Kann, A. Marchetti-Spaccamela, and M. Protasi, Complexity and Approximation (Springer, Berlin, 1999).
- [2] M. Baroni, C. Semple, and M. Steel, A framework for representing reticulate evolution, Ann. Combin., in press.
- [3] M. Baroni, S. Grünewald, V. Moulton, and C. Semple, Bounding the number of hybridisation events for a consistent evolutionary history, submitted.
- [4] M. Bordewich and C. Semple, On the computational complexity of the rooted subtree prune and regraft distance, Ann. Combin., in press.

- [5] M. Chlebík and J. Chlebková, Inapproximability results for bounded variants of optimization problems, in: A. Lingas and B. J. Nilsson, eds, *Fundamentals of Computation Theory, 14th International Symposium (FCT)*, Lecture Notes in Computer Science, Vol. 2751 (Springer-Verlag, 2003) 27-38.
- [6] S. Grünwald, Private communication.
- [7] D. Gusfield, S. Eddhu, and C. Langley, Optimal, efficient reconstruction of phylogenetic networks with constrained recombination, *J. Bioinformatics and Computational Biology* 2 (2004) 173-213.
- [8] J. Hein, Reconstructing evolution of sequences subject to recombination using parsimony, *Math. Biosci.* 98 (1990) 185-200.
- [9] J. Hein, T. Jing, L. Wang, K. Zhang, On the complexity of comparing evolutionary trees, *Discrete Appl. Math.* 71 (1996) 153-169.
- [10] V. Kann, Maximum bounded 3-dimensional matching is MAX SNP-complete, *Inform. Process. Lett.* 37 (1991) 27-35.
- [11] W. Maddison, Gene trees in species trees, *Syst. Biol.* 46 (1997) 523-536.
- [12] L. Nakhleh, T. Warnow, and C. Randal Linder, Reconstructing reticulate evolution in species - theory and practice, in: *Proceedings of the 8th Annual International Conference on Research in Computational Molecular Biology (RECOMB)* (2004) 337-346.
- [13] C. H. Papadimitriou and M. Yannakakis, Optimization, approximation, and complexity classes, *J. Comput. System Sci.* 43 (1991) 425-440.
- [14] E. M. Rodrigues, M.-F. Sagot, and Y. Wakabayashi, Some approximation results for the maximum agreement forest problem, in: M. Goemans *et al.*, eds, *Approximation, Randomization and Combinatorial Optimization: Algorithms and Techniques (APPROX and RANDOM)*, Lecture Notes in Computer Science, Vol. 2129 (Springer, Berlin, 2001) 159-169.
- [15] C. Semple and M. Steel, *Phylogenetics* (Oxford University Press, 2003).
- [16] Y. Song and J. Hein, Parsimonious reconstruction of sequence evolution and haplotype blocks: finding the minimum number of recombination events, in: G. Benson and R. Page, eds, *Algorithms in Bioinformatics (WABI)*, Lecture Notes in Bioinformatics, Vol. 2812 (Springer, Berlin, 2003) 287-302.
- [17] Y. Song and J. Hein, Constructing minimal ancestral recombination graphs, *J. Comp. Biol.*, in press.
- [18] L. Wang, K. Zhang, and L. Zhang, Perfect phylogenetic networks with recombination, *J. Computational Biology* 8 (2001) 69-78.

BIOMATHEMATICS RESEARCH CENTRE, DEPARTMENT OF MATHEMATICS AND STATISTICS, UNIVERSITY OF CANTERBURY, CHRISTCHURCH, NEW ZEALAND

Current address: School of Computing, University of Leeds, Leeds LS2 9JT, UK

E-mail address: magnusb@comp.leeds.ac.uk

BIOMATHEMATICS RESEARCH CENTRE, DEPARTMENT OF MATHEMATICS AND STATISTICS, UNIVERSITY OF CANTERBURY, CHRISTCHURCH, NEW ZEALAND

E-mail address: c.semple@math.canterbury.ac.nz